Optimization for Machine Learning "Hands On"

Alexandre Gramfort (Inria) http://alexandre.gramfort.net/

Quentin Bertrand (Inria) https://qb3.github.io

Disclaimer

- 3 hours is too short for a detailed course on optimization
- We will cover only unconstrained problems
- ▶ We will not cover non-smooth problems (e.g., Lasso, SVM)
- The objective is to grasp quickly some theoretical aspects ...
- ... and to code everything to be able to experiment.

Introduction

Gradient descent

Newton method

Stochastic gradient descent

Coordinate descent

References and useful links

- S. Boyd and L. Vandenberghe. Convex optimization. Cambridge University Press, 2004, pp. xiv+716
- J. Nocedal and S. J. Wright. Numerical optimization.
 Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006
- Cool website for visualization of optimization algorithms: http://fa.bianp.net/teaching/2018/COMP-652/

Is this a cat?





Is this a cat?



Yes

Is this a cat?



x: Input / Feature y: Output / Target

Goal: find mapping h that assigns the 'correct' target to each input

$$h: x \in \mathbb{R}^p \longrightarrow y \in \mathbb{R}$$

Empirical Risk Minimization (ERM)

Goal: from examples $(x_1, y_1), \ldots, (x_n, y_n)$ learn a function $h : \mathbb{R}^p \to \mathbb{R}$ such that

 $h(x_{n+1}) \simeq y_{n+1}$

Ideally, for a given loss function L :

$$h^* \in \underset{h \in \mathcal{H}}{\arg\min} \underbrace{\mathbb{E}[L(h(x), y)]}_{\text{Expected risk}}$$

Empirical Risk Minimization (ERM)

Goal: from examples $(x_1, y_1), \ldots, (x_n, y_n)$ learn a function $h : \mathbb{R}^p \to \mathbb{R}$ such that

 $h(x_{n+1}) \simeq y_{n+1}$

Ideally, for a given loss function L :

$$h^* \in \underset{h \in \mathcal{H}}{\operatorname{arg\,min}} \underbrace{\mathbb{E}[L(h(x), y)]}_{\text{Expected risk}}$$

In practice:

$$h^* \in \underset{h \in \mathcal{H}}{\operatorname{arg\,min}} \underbrace{\frac{1}{n} \sum_{i=1}^{n} L(h(x_i), y_i)}_{\operatorname{Empirical risk}}$$

Empirical Risk Minimization (ERM)

Goal: from examples $(x_1, y_1), \ldots, (x_n, y_n)$ learn a function $h : \mathbb{R}^p \to \mathbb{R}$ such that

 $h(x_{n+1}) \simeq y_{n+1}$

Ideally, for a given loss function L:

$$h^* \in \underset{h \in \mathcal{H}}{\operatorname{arg\,min}} \underbrace{\mathbb{E}[L(h(x), y)]}_{\text{Expected risk}}$$

In practice:

$$h^* \in \underset{h \in \mathcal{H}}{\operatorname{arg\,min}} \underbrace{\frac{1}{n} \sum_{i}^{n} L(h(x_i), y_i)}_{\text{Empirical risk}}$$

Examples of ERM in Practice

Let $X = [x_1, \dots, x_n]^\top \in \mathbb{R}^{n \times p}$ (design matrix) (Regularized) Linear Regression: $w^* = \underset{w \in \mathbb{R}^p}{\operatorname{arg\,min}} \underbrace{\frac{1}{2} \|y - Xw\|^2}_{\ell_2 \operatorname{Loss}} + \frac{\lambda}{2} \|w\|^2$

(Regularized) Logistic Regression:

$$w^* = \underset{w \in \mathbb{R}^p}{\operatorname{arg\,min}} \sum_{i=1}^{n} \underbrace{\log(1 + \exp(-y_i x_i^{\top} w))}_{\operatorname{Logistic\,Loss}} + \frac{\lambda}{2} \|w\|^2$$

Examples of ERM in Practice

Let $X = [x_1, \dots, x_n]^\top \in \mathbb{R}^{n \times p}$ (design matrix)

(Regularized) Linear Regression:

$$w^* = \underset{w \in \mathbb{R}^p}{\operatorname{arg\,min}} \underbrace{\frac{1}{2} \|y - Xw\|^2}_{\ell_2 \operatorname{Loss}} + \frac{\lambda}{2} \|w\|^2$$

(Regularized) Logistic Regression:

$$w^* = \underset{w \in \mathbb{R}^p}{\operatorname{arg\,min}} \sum_{i=1}^{n} \underbrace{\log(1 + \exp(-y_i x_i^{\top} w))}_{\operatorname{Logistic\,Loss}} + \frac{\lambda}{2} \|w\|^2$$

All of these are convex optimization problems!

Examples of ERM in Practice

Let $X = [x_1, \dots, x_n]^\top \in \mathbb{R}^{n \times p}$ (design matrix)

(Regularized) Linear Regression:

$$w^* = \underset{w \in \mathbb{R}^p}{\operatorname{arg\,min}} \underbrace{\frac{1}{2} \|y - Xw\|^2}_{\ell_2 \operatorname{Loss}} + \frac{\lambda}{2} \|w\|^2$$

(Regularized) Logistic Regression:

$$w^* = \underset{w \in \mathbb{R}^p}{\operatorname{arg\,min}} \sum_{i=1}^{n} \underbrace{\log(1 + \exp(-y_i x_i^{\top} w))}_{\operatorname{Logistic\,Loss}} + \frac{\lambda}{2} \|w\|^2$$

All of these are convex optimization problems!

Gradient

Definition (Gradient) For $f : \mathbb{R}^p \to \mathbb{R}$ smooth the gradient reads: $\nabla f(w) = \left[\frac{\partial f(w)}{\partial w_1}, \dots, \frac{\partial f(w)}{\partial w_p}\right]^\top \in \mathbb{R}^p$

Examples

•
$$f(w) = x^{\top}w$$
 where $x \in \mathbb{R}^p$, then $\nabla f(w) = x$

•
$$f(w) = g(x^{\top}w)$$
 where $g : \mathbb{R} \to \mathbb{R}$, then $\nabla f(w) = g'(x^{\top}w)x$

•
$$f(w) = w^{\top}Aw$$
 where $A \in \mathbb{R}^{p \times p}$, $\nabla f(w) = (A + A^{\top})w$

Hessian



Examples

•
$$f(w) = x^{\top}w$$
 where $x \in \mathbb{R}^p$, then $\nabla^2 f(w) = 0$

•
$$f(w) = w^{\top}Aw$$
 where $A \in \mathbb{R}^{p \times p}$, $\nabla^2 f(w) = A + A^{\top}$

Warm up!

You have 3 minutes to compute the gradient ∇f and the Hessian ∇²f of the linear regression function:

$$f: w \mapsto \frac{1}{2} \|y - Xw\|^2$$

Hint:

$$f(w) = \frac{1}{2} (y - Xw)^{\top} (y - Xw)$$

= $\frac{1}{2} ||y||^2 - (X^{\top}y)^{\top}w + \frac{1}{2}w^{\top}X^{\top}Xw$

Warm up!

You have 3 minutes to compute the gradient ∇f and the Hessian ∇²f of the linear regression function:

$$f: w \mapsto \frac{1}{2} \|y - Xw\|^2$$

Hint:

$$f(w) = \frac{1}{2}(y - Xw)^{\top}(y - Xw)$$

= $\frac{1}{2}||y||^2 - (X^{\top}y)^{\top}w + \frac{1}{2}w^{\top}X^{\top}Xw$

Solution:

$$\nabla f(w) = X^{\top} (Xw - y)$$
$$\nabla^2 f(w) = X^{\top} X$$

Convexity I

Definition (Convex function)

$$\begin{split} f: \mathbb{R}^p \to \mathbb{R} \text{ is convex if and only if, } \forall u, v \in \mathbb{R}^p \text{, } \forall \lambda \in [0,1]\text{:} \\ f(\lambda u + (1-\lambda)v) \leq \lambda f(u) + (1-\lambda)f(v) \end{split}$$



Convexity II

Proposition (Convex differentiable function / has a gradient)

$$\begin{split} f: \mathbb{R}^p \to \mathbb{R} \text{ is convex if and only if, } \forall u, v \in \mathbb{R}^p: \\ f(v) \geq f(u) + \nabla f(u)^\top (v-u) \end{split}$$



Convexity III

Proposition (Convex twice differentiable function)

$f: \mathbb{R}^p \to \mathbb{R}$ is convex if and only if, $\forall u \in \mathbb{R}^p$: $\nabla^2 f(u) \succeq 0$



Strong Convexity

Definition (Strongly convex function)

$$\begin{split} f: \mathbb{R}^p &\to \mathbb{R} \text{ is } \mu\text{-strongly convex if and only if, } \forall u, v \in \mathbb{R}^p: \\ f(v) &\geq f(u) + \nabla f(u)^\top (v-u) + \frac{\mu}{2} \|v-u\|^2 \end{split}$$



Smoothness

Definition (L-Smoothness)

 $\begin{array}{l} \mathsf{A} \text{ function } f: \mathbb{R}^p \to \mathbb{R} \text{ is } L\text{-smooth if } \nabla f \text{ is } L\text{-Lipschitz:} \\ \|\nabla f(u) - \nabla f(v)\| \leq L \left\|u - v\right\|, \quad \forall u, v \in \mathbb{R}^p \end{array}$

in particular this implies

$$f(v) \le f(u) + \nabla f(u)^{\top}(v-u) + \frac{L}{2} ||v-u||^2$$

Remark: In practice one wants L as small as possible



L-Smoothness & Strong Convexity & Hessian

Proposition (L-Smoothness twice differentiable function)

 $f: \mathbb{R}^p \to \mathbb{R}$ is L-smooth convex if and only if, $\forall u \in \mathbb{R}^p: \nabla^2 f(u) \preceq L \operatorname{Id}_p$

Proposition (Strongly convex twice differentiable function) $f : \mathbb{R}^p \to \mathbb{R}$ is μ -strongly convex if and only if, $\forall u \in \mathbb{R}^p$: $\nabla^2 f(u) \succeq \mu \operatorname{Id}_p$



From surrogate minimization to Gradient Descent (GD)

For f L-smooth Taylor expansion gives:

$$f(w) \le f(w^k) + \nabla f(w^k)^\top (w - w^k) + \frac{L}{2} ||w - w^k||^2$$

Surrogate function $\tilde{f}(w)$

Idea: Minimize the surrogate function \tilde{f}

 \overline{f} is convex, differentiable, infinite at the infinite (coercive) \Rightarrow its minimum w^* is achieved where gradient is 0:

$$0 = \nabla \tilde{f}(w^*) = \nabla f(w^k) + L(w^* - w^k)$$

From surrogate minimization to Gradient Descent (GD)

For f L-smooth Taylor expansion gives:

$$f(w) \le \underbrace{f(w^k) + \nabla f(w^k)^\top (w - w^k) + \frac{L}{2} \|w - w^k\|^2}_{=}$$

Surrogate function $\tilde{f}(w)$

Idea: Minimize the surrogate function \tilde{f}

 \tilde{f} is convex, differentiable, infinite at the infinite (coercive) \Rightarrow its minimum w^* is achieved where gradient is 0:

$$0 = \nabla \tilde{f}(w^*) = \nabla f(w^k) + L(w^* - w^k)$$

Taking w^{k+1} as the minimizer of $ar{f}$ leads to gradient descent:

$$w^{k+1} = w^k - \frac{1}{L}\nabla f(w^k)$$

From surrogate minimization to Gradient Descent (GD)

For f L-smooth Taylor expansion gives:

$$f(w) \le \underbrace{f(w^k) + \nabla f(w^k)^\top (w - w^k) + \frac{L}{2} \|w - w^k\|^2}_{=}$$

Surrogate function $\tilde{f}(w)$

Idea: Minimize the surrogate function \tilde{f}

 \tilde{f} is convex, differentiable, infinite at the infinite (coercive) \Rightarrow its minimum w^* is achieved where gradient is 0:

$$0 = \nabla \tilde{f}(w^*) = \nabla f(w^k) + L(w^* - w^k)$$

Taking w^{k+1} as the minimizer of \tilde{f} leads to gradient descent:

$$w^{k+1} = w^k - \frac{1}{L}\nabla f(w^k)$$

Exercise 1

Exercise Write the GD algorithm for the linear regression

$$f: w \mapsto \frac{1}{2} \|y - Xw\|^2$$

Algorithm: Gradient Descent

 $\begin{array}{ll} \text{init} & : w^0 = 0_p, \ L \\ \text{for iter} = 1, \dots, \ \text{do} \\ \mid & w^{k+1} = w^k - \frac{1}{L} \nabla f(w^k) \end{array}$

Exercise 1

Solution Write the GD algorithm for the linear regression

$$f: w \mapsto \frac{1}{2} \|y - Xw\|^2$$

Algorithm: Gradient Descent

init :
$$w^0 = 0_p$$
, $L := \lambda_{\max}(X^\top X) = ||X||_2^2$
for iter = 1,..., do
 $| w^{k+1} = w^k - \frac{1}{L}X^\top (Xw^k - y)$

Line-search⁽¹⁾

What to do when a function is smooth, but you do not know the Lipschitz constant L? Or when L is too conservative?

Gradient descent with variable stepsize:

$$w^{k+1} = w^k - \alpha^k \nabla f(w^k)$$

where the stepsize α^k changes at each iteration and is found by line-search.

⁽¹⁾ J. Nocedal and S. J. Wright. *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006, Chap. 3.

Line-search⁽¹⁾

What to do when a function is smooth, but you do not know the Lipschitz constant L? Or when L is too conservative?

Gradient descent with variable stepsize:

$$w^{k+1} = w^k - \alpha^k \nabla f(w^k)$$

where the stepsize α^k changes at each iteration and is found by line-search.

⁽¹⁾ J. Nocedal and S. J. Wright. *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006, Chap. 3.

Hands on 0

 \rightarrow See notebook: 00-gradient_descent_line_search.ipynb

Gradient descent: theoretical results

Algorithm: GD

init :
$$w^0 = 0_p$$
, L
for iter = 1, ..., do
 $| w^{k+1} = w^k - \frac{1}{L} \nabla f(w^k)$

Proposition

If *f* is *µ*-strongly convex and *L*-smooth, then: $\| k \| \leq \left(1 - \frac{\mu}{k} \right)^k \| = 0$

$$|w^k - w^*|| \le \left(1 - \frac{\mu}{L}\right)^* ||w^0 - w^*|$$

One has "linear" (a.k.a. "exponential") convergence.

Proposition

If f is convex and L-smooth, then:
$$f(w^k) - f(w^*) \leq \frac{2L \|w^0 - w^*\|^2}{k}$$

One has "sublinear" convergence.

Gradient descent: theoretical results

Algorithm: GD

init :
$$w^0 = 0_p$$
, L
for iter = 1,..., do
 $| w^{k+1} = w^k - \frac{1}{L} \nabla f(w^k)$

Proposition

If f is μ -strongly convex and L-smooth, then:

$$||w^k - w^*|| \le \left(1 - \frac{\mu}{L}\right)^{\kappa} ||w^0 - w^*||$$

One has "linear" (a.k.a. "exponential") convergence.

Proposition

If f is convex and L-smooth, then: $f(w^k) - f(w^*) \leq \frac{2L \left\|w^0 - w^*\right\|^2}{k}$

One has "sublinear" convergence.

Hands on 1

 \rightarrow See notebook: O1-logistic_gd.ipynb

Remark

- ▶ In all the hands on the main algorithm is already implemented
- We propose you to add little modifications
- ► All the solutions are in the solutions folder
- But do not look at them too quickly ...

Newton method: intuition

Taylor expansion to the order 2:

$$f(w) \approx \underbrace{f(w^k) + \nabla f(w^k)^\top (w - w^k) + \frac{1}{2} (w - w^k)^\top \nabla^2 f(w^k) (w - w^k)}_{\tilde{f}(w)}$$

Idea: Minimize the quadratic approximation f:

$$0 = \nabla \tilde{f}(w^{k+1}) = \nabla f(w^k) + \nabla^2 f(w^k)(w^{k+1} - w^k)$$

Newton method: intuition

Taylor expansion to the order 2:

$$f(w) \approx \underbrace{f(w^{k}) + \nabla f(w^{k})^{\top}(w - w^{k}) + \frac{1}{2}(w - w^{k})^{\top} \nabla^{2} f(w^{k})(w - w^{k})}_{\tilde{f}(w)}$$

Idea: Minimize the quadratic approximation \tilde{f} :

$$0=\nabla \tilde{f}(w^{k+1})=\nabla f(w^k)+\nabla^2 f(w^k)(w^{k+1}-w^k)$$

This leads to the following update for the Newton algorithm:

$$w^{k+1} = w^k - [\nabla^2 f(w^k)]^{-1} \nabla f(w^k)$$
Newton method: intuition

Taylor expansion to the order 2:

$$f(w) \approx \underbrace{f(w^k) + \nabla f(w^k)^\top (w - w^k) + \frac{1}{2} (w - w^k)^\top \nabla^2 f(w^k) (w - w^k)}_{\tilde{f}(w)}$$

Idea: Minimize the quadratic approximation \tilde{f} :

$$0 = \nabla \tilde{f}(w^{k+1}) = \nabla f(w^k) + \nabla^2 f(w^k)(w^{k+1} - w^k)$$

This leads to the following update for the Newton algorithm:

$$w^{k+1} = w^k - [\nabla^2 f(w^k)]^{-1} \nabla f(w^k)$$

What about convergence guarantees?

Newton method: intuition

Taylor expansion to the order 2:

$$f(w) \approx \underbrace{f(w^k) + \nabla f(w^k)^\top (w - w^k) + \frac{1}{2} (w - w^k)^\top \nabla^2 f(w^k) (w - w^k)}_{\tilde{f}(w)}$$

Idea: Minimize the quadratic approximation \tilde{f} :

$$0 = \nabla \tilde{f}(w^{k+1}) = \nabla f(w^k) + \nabla^2 f(w^k)(w^{k+1} - w^k)$$

This leads to the following update for the Newton algorithm:

$$w^{k+1} = w^k - [\nabla^2 f(w^k)]^{-1} \nabla f(w^k)$$

What about convergence guarantees?

Exercise Write the Newton algorithm for the linear regression

$$f:w \ \mapsto \frac{1}{2}\|y-Xw\|^2$$

Algorithm: Newton algorithm

 $\begin{array}{ll} \text{init} & : w^0 = 0_p \\ \text{for iter} = 1, \dots, \text{ do} \\ \mid & w^{k+1} = w^k - [\nabla^2 f(w^k)]^{-1} \nabla f(w^k) \end{array}$

Solution Write the Newton algorithm for the linear regression

$$f:w \ \mapsto \frac{1}{2}\|y-Xw\|^2$$

Algorithm: Newton algorithm

 $\begin{array}{ll} \text{init} & : w^0 = 0_p \\ \text{for iter} = 1, \dots, \text{ do} \\ & \mid \ w^{k+1} = w^k - [X^\top X]^{-1} X^\top (X w^k - y) \end{array}$

Solution Write the Newton algorithm for the linear regression

$$f: w \mapsto \frac{1}{2} \|y - Xw\|^2$$

Algorithm: Newton algorithm

init :
$$w^0 = 0_p$$

for iter = 1,..., do
 $| w^{k+1} = w^k - [X^T X]^{-1} X^T (X w^k - y)$

Question: Is there an interest of applying Newton method on a linear regression problem?

Convergence of Newton method

Theorem (Convergence of Newton method)

Suppose that f is twice differentiable, and the Hessian $\nabla^2 f(w^*)$ is Lipschitz continuous in a neighborhood of a solution w^* such that $\nabla f(w^*) = 0$ and $\nabla^2 f(w^*) \succ 0$. Then there exists a closed ball \mathcal{B} centered on w^* , such that for every $w^0 \in \mathcal{B}$, the sequence w^k obtained with Newton algorithm stays in \mathcal{B} and converges towards w^* . Furthermore, there is a constant $\gamma > 0$, such that $||w^{k+1} - w^*|| \le \gamma ||w^k - w^*||^2$. One has "super linear" convergence.

Drawbacks:

- ► Convergence of Newton is local (see proof in⁽²⁾). The method may diverge if the initial point is too far from w^{*} or if the Hessian is not positive definite
- One has to solve a linear system at each step! $O(p^3)$

⁽²⁾J. Nocedal and S. J. Wright. *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006, Thm. 3.5.

Convergence of Newton method

Theorem (Convergence of Newton method)

Suppose that f is twice differentiable, and the Hessian $\nabla^2 f(w^*)$ is Lipschitz continuous in a neighborhood of a solution w^* such that $\nabla f(w^*) = 0$ and $\nabla^2 f(w^*) \succ 0$. Then there exists a closed ball \mathcal{B} centered on w^* , such that for every $w^0 \in \mathcal{B}$, the sequence w^k obtained with Newton algorithm stays in \mathcal{B} and converges towards w^* . Furthermore, there is a constant $\gamma > 0$, such that $||w^{k+1} - w^*|| \le \gamma ||w^k - w^*||^2$. One has "super linear" convergence.

Drawbacks:

- Convergence of Newton is local (see proof in⁽²⁾). The method may diverge if the initial point is too far from w* or if the Hessian is not positive definite
- One has to solve a linear system at each step! $O(p^3)$

⁽²⁾ J. Nocedal and S. J. Wright. *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006, Thm. 3.5.

From Newton to quasi-Newton

Idea:

- Construct iterative approximations of (the inverse of) the Hessian
- Combine it with line-search strategies

 $\begin{cases} d^k &= -B^k \nabla f(w^k) & \text{find a descent direction} \\ w^{k+1} &= w^k + \alpha^k d^k & \text{line-search} \end{cases}$

⁽³⁾ J. Nocedal and S. J. Wright. Numerical optimization. Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006.

⁽⁴⁾ J. Nocedal. "Updating quasi-Newton matrices with limited storage". In: Mathematics of computation 35.151 (1980), pp. 773–782.

From Newton to quasi-Newton

Idea:

- Construct iterative approximations of (the inverse of) the Hessian
- Combine it with line-search strategies

$$\begin{cases} d^k &= -B^k \nabla f(w^k) & \text{find a descent direction} \\ w^{k+1} &= w^k + \alpha^k d^k & \text{line-search} \end{cases}$$

- ▶ There exists a whole jungle of iterative approximations for the inverse of the Hessian: $B^{k+1} = B^k + \Delta^{k}$ ⁽³⁾
- The one you should know about is BFGS strategy and the memory efficient L-BFGS⁽⁴⁾

^{(&}lt;sup>3)</sup> J. Nocedal and S. J. Wright. *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006.

⁽⁴⁾ J. Nocedal. "Updating quasi-Newton matrices with limited storage". In: Mathematics of computation 35.151 (1980), pp. 773–782.

From Newton to quasi-Newton

Idea:

- Construct iterative approximations of (the inverse of) the Hessian
- Combine it with line-search strategies

$$\begin{cases} d^k &= -B^k \nabla f(w^k) & \text{find a descent direction} \\ w^{k+1} &= w^k + \alpha^k d^k & \text{line-search} \end{cases}$$

- ► There exists a whole jungle of iterative approximations for the inverse of the Hessian: $B^{k+1} = B^k + \Delta^{k}$ ⁽³⁾
- The one you should know about is BFGS strategy and the memory efficient L-BFGS⁽⁴⁾

⁽³⁾ J. Nocedal and S. J. Wright. *Numerical optimization*. Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006.

⁽⁴⁾ J. Nocedal. "Updating quasi-Newton matrices with limited storage". In: Mathematics of computation 35.151 (1980), pp. 773–782.

Hands on 2

 \rightarrow See notebook: 02-logistic_newton.ipynb

Remark

► All the solutions are in the solutions folder ...

Optimization for ML: finite sum

Optimization in general:

 $\min_{x\in\mathcal{H}}f(x)$

Optimization for Machine Learning (linear or logistic regression): $\min_{x\in\mathbb{R}^p}f(x):=\sum_{i=1}^n f_i(x)$

Optimization for ML: finite sum

Optimization in general:

 $\min_{x \in \mathcal{H}} f(x)$

Optimization for Machine Learning (linear or logistic regression): $\min_{x\in\mathbb{R}^p}f(x):=\sum_{i=1}^n f_i(x)$

Can we take advantage of this finite sum structure?

▶ In particular when the number of samples *n* is large?

Optimization for ML: finite sum

Optimization in general:

 $\min_{x \in \mathcal{H}} f(x)$

Optimization for Machine Learning (linear or logistic regression): $\min_{x\in\mathbb{R}^p}f(x):=\sum_{i=1}^nf_i(x)$

- Can we take advantage of this finite sum structure?
- ▶ In particular when the number of samples n is large?

Stochastic Gradient Descent (SGD)

Question: what is the cost per iteration of one update of GD?

- Full gradient methods can be expensive O(np) per iteration!
- ▶ Idea: use a single gradient $\nabla f_i(w)$ instead of full gradient $\sum_{i=1}^{n} \nabla f_i(w)$ (*n* times faster!)

Algorithm: SGDinit : $w = 0_p$, α for iter = 1, ..., doSample i in 1, ..., n// Single gradient $\nabla f_i(w)$ call// O(p) per update $w \leftarrow w - \alpha \nabla f_i(w)$ return w

Algorithm: GD init : $w = 0_p$, α for iter = 1,..., do // Full gradient $\nabla f(w)$ call // O(np) per update $w \leftarrow w - \alpha \nabla f(w)$ return w

Stochastic Gradient Descent (SGD)

Question: what is the cost per iteration of one update of GD?

- Full gradient methods can be expensive O(np) per iteration!
- ▶ Idea: use a single gradient $\nabla f_i(w)$ instead of full gradient $\sum_i^n \nabla f_i(w)$ (*n* times faster!)

Algorithm: SGD

Algorithm: GDinit: $w = 0_p, \alpha$ for iter = 1, ..., do// Full gradient $\nabla f(w)$ call// O(np) per update $w \leftarrow w - \alpha \nabla f(w)$ return w

Exercise Write the algorithm for the linear regression

$$f: w \mapsto \sum_{i=1}^{n} \underbrace{\frac{1}{2} (y_i - x_i^{\top} w)^2}_{f_i(w)}$$

Algorithm: SGD

init :
$$w = 0_p$$
, α
for iter = 1, ..., do
Sample *i* in 1, ..., *n*
// Single gradient $\nabla f_i(w)$ call
// $O(p)$ per update
 $w \leftarrow w - \alpha \nabla f_i(w)$
return w

Algorithm: GD

$$\begin{array}{ll} \text{for iter} & : & w = 0_p, \ \alpha \\ \text{for iter} & = 1, \dots, \ \text{do} \\ \\ & & \swarrow^{// \ \text{Full gradient } \nabla f(w) \ \text{call}} \\ & & \swarrow^{O(np) \ \text{per update}} \\ & & w \leftarrow w - \alpha \nabla f(w) \\ \text{return } w \end{array}$$

Solution Write the algorithm for the linear regression

$$f: w \mapsto \sum_{i=1}^{n} \underbrace{\frac{1}{2} (y_i - x_i^{\top} w)^2}_{f_i(w)}$$

Algorithm: SGD

Algorithm: GD

init :
$$w = 0_p$$
, $\alpha = 1/||X||_2^2$
for iter = 1,..., do
// Full gradient $\nabla f(w)$ call
// $O(np)$ per update
 $w \leftarrow w - \alpha X^{\top}(Xw - y)$
return w

Solution Write the algorithm for the linear regression

$$f: w \mapsto \sum_{i=1}^{n} \underbrace{\frac{1}{2} (y_i - x_i^{\top} w)^2}_{f_i(w)}$$

Algorithm: SGD

Algorithm: GD init : $w = 0_p$, $\alpha = 1/||X||_2^2$ for iter = 1, ..., do // Full gradient $\nabla f(w)$ call // O(np) per update $w \leftarrow w - \alpha X^{\top}(Xw - y)$ return w

Question: how to choose α for the SGD?

SGD with constant stepsize $\alpha {:}\ {\rm theory}$

Proposition

If f is μ -strongly convex and L-smooth, then: $\mathbb{E}(\|w^k - w^*\|^2) \le (1 - \alpha \mu)^k \|w^0 - w^*\|^2 + \frac{\alpha}{\mu}C$

Remark

- ► Constant C depends on the norm for the f_i gradients (bounded gradient hypothesis)
- SGD with constant step size α does not converge
- In ML if the estimation and the approximation are bigger than the optimization error it's ok!^a

^aL. Bottou and O. Bousquet. "The tradeoffs of large scale learning". In: Advances in neural information processing systems. 2008, pp. 161–168.

Example of SGD



Example of SGD



Example of SGD



Hands on 3

 \rightarrow See notebook: 03-stochastic_gradient_descent.ipynb

Remark

► All the solutions are in the solutions folder ...

Goal:

 $\min_{w\in\mathbb{R}^p}f(w_1,\ldots,w_p)$

Idea: solve smaller and simpler problems (one coordinate at a time)

Goal:

$$\min_{w\in\mathbb{R}^p}f(w_1,\ldots,w_p)$$

Idea: solve smaller and simpler problems (one coordinate at a time)

Algorithm: Exact coordinate descent init : $w = 0_p$ for iter = 1,..., do for j = 1, ..., p do $\bigvee_{j \in arg \min_{z \in \mathbb{R}}} f(w_1, ..., w_{j-1}, z, w_{j+1}, ..., w_p)$ return w

Goal:

$$\min_{w\in\mathbb{R}^p}f(w_1,\ldots,w_p)$$

Idea: solve smaller and simpler problems (one coordinate at a time)

Algorithm: Exact coordinate descent

 $\begin{array}{ll} \text{init} & : w = 0_p \\ \text{for iter} = 1, \dots, \text{ do} \\ & & \left| \begin{array}{c} \text{for } j = 1, \dots, p \text{ do} \\ & & | \end{array} \right|^{// \text{Update only one coef.}} \\ & & w_j \leftarrow \arg\min_{z \in \mathbb{R}} f(w_1, \dots, w_{j-1}, z, w_{j+1}, \dots, w_p) \\ \text{return } w \end{array}$

Remark

- ▶ The order of cycle through coordinates is arbitrary, can use any permutation of 1, 2, ..., p
- We just have to solve 1D optimization problems but a lot of them...

Goal:

$$\min_{w\in\mathbb{R}^p}f(w_1,\ldots,w_p)$$

Idea: solve smaller and simpler problems (one coordinate at a time)

Algorithm: Exact coordinate descent

 $\begin{array}{ll} \text{init} & : w = 0_p \\ \text{for iter} = 1, \dots, \text{ do} \\ & & \left| \begin{array}{c} \text{for } j = 1, \dots, p \text{ do} \\ & & | \end{array} \right|^{// \text{Update only one coef.}} \\ & & w_j \leftarrow \arg\min_{z \in \mathbb{R}} f(w_1, \dots, w_{j-1}, z, w_{j+1}, \dots, w_p) \\ \text{return } w \end{array}$

Remark

- \blacktriangleright The order of cycle through coordinates is arbitrary, can use any permutation of $1,2,\ldots,p$
- We just have to solve 1D optimization problems but a lot of them...

Example of CD



Coordinate Gradient Descent: algorithm

- Exact minimization can be expensive
- Idea: do a local gradient step instead of exact minimization
- Step-sizes are now given by the coordinate-wise functions

Algorithm: CD	Algorithm: GD
init : $w = 0_p, L_1, \ldots, L_p$	init : $w = 0_p$, L
for $iter = 1, \ldots, do$	for $iter = 1, \ldots, do$
for $j = 1, \dots, p$ do	// Update all the coef. $w \leftarrow w - \frac{1}{T} \nabla f(w)$
$w_j \leftarrow w_j - \frac{1}{L_j} \nabla_j f(w)$	return w
return w	

▶ Exercise Write the CD algorithm for the linear regression

$$f: w \mapsto \|y - Xw\|^2$$

Algorithm: CD	Algorithm: GD
init : $w = 0_p, L_1, \dots, L_p$	init : $w = 0_p$, L
for iter = 1,, do	for iter = 1,, do
for $j = 1, \dots, p$ do	$ // Full gradient \nabla f(w) call$
// Update only one coef.	$w \leftarrow w - \frac{1}{L} \nabla f(w)$
$w_j \leftarrow w_j - \frac{1}{L_j} \nabla_j f(w)$	return w

Solution Write the CD algorithm for the linear regression

$$f: w \mapsto \|y - Xw\|^2$$

Algorithm: CD

init :
$$w = 0_p, L_j = ||X_{:j}||^2$$

for iter = 1,..., do
for $j = 1, ..., p$ do
 $| // Update only one coef.$
 $w_j \leftarrow w_j - \frac{X_{:j}^\top(Xw-y)}{L_j}$

Algorithm: GD

 $\begin{array}{ll} \text{init} & : \ w = 0_p, \ \alpha = 1/\|X\|_2^2 \\ \text{for iter} = 1, \dots, \ \text{do} \\ & | & \\ // \ \text{Full gradient} \ \nabla f(w) \ \text{call} \\ & w \leftarrow w - \alpha X^\top (Xw - y) \\ \text{return } w \end{array}$

Solution Write the CD algorithm for the linear regression

$$f: w \mapsto \|y - Xw\|^2$$

Algorithm: CD

$$\begin{array}{ll} \text{init} & : w = 0_p, \ L_j = \|X_{:j}\|^2 \\ \text{for iter} = 1, \dots, \ \textbf{do} \\ & & \\ & \text{for } j = 1, \dots, p \ \textbf{do} \\ & & \\ &$$

Algorithm: GD

 $\begin{array}{ll} \text{init} & : \ w = 0_p, \ \alpha = 1/\|X\|_2^2 \\ \text{for iter} = 1, \ldots, \ \text{do} \\ & & \\ & | \begin{array}{c} \ \prime / \ \text{Full gradient} \ \nabla f(w) \ \text{call} \\ & w \leftarrow w - \alpha X^\top (Xw - y) \\ \text{return } w \end{array} \end{array}$

Question What is the cost of one update of CD?

Convergence speed of CD⁽⁶⁾

Assume f is convex; ∇f is Lipschitz continuous

Proposition (Beck and Tetruashvili (2013))

$$f(w^{k+1}) - f(w^*) \le 4L_{\max}(1 + n^3 L_{\max}^2 / L_{\min}^2) \frac{R^2(w^0)}{k + 8/n}$$

where $R^2(w^0) = \max_{u,v \in \mathcal{V}} \{ \|u - v\| : f(v) \le f(u) \le f(w^0) \}$,
 $L_{\max} = \max_j L_j$ and $L_{\min} = \min_j L_j$.

- ► Worst case complexity rates of CD are bad (n³ complexity) because it is possible to construct adversarial examples⁽⁵⁾
- Yet CD performs surprisingly well on real-world problems (see hands on)
- ▶ Random coordinate selection has a better average complexity

⁽⁵⁾R. Sun and Y. Ye. "Worst-case complexity of cyclic coordinate descent: $O(n^2)$ gap with randomized version". In: *Mathematical Programming* (2019), pp. 1–34.

⁽⁶⁾A. Beck and L. Tetruashvili. "On the convergence of block coordinate type methods". In: SIAM J. Imaging Sci. 23.4 (2013), pp. 651–694.

Take a look back before "Hands on 4"

$$f: w \mapsto \frac{1}{2} \|y - Xw\|^2$$

Algorithm: GD

 $\begin{array}{ll} \text{init} & : \ w = 0_p, \ \alpha = 1/\|X\|_2^2 \\ \text{for iter} = 1, \dots, \ \text{do} \\ & | \ // \ \text{Full gradient} \ \nabla f(w) \ \text{call} \\ & w \leftarrow w - \alpha X^\top (Xw - y) \\ \text{return} \ w \end{array}$

Algorithm: SGD

Algorithm: CD

return w

Hands on 4

 \rightarrow See notebook: 04-coordinate_descent.ipynb

Remark

► All the solutions are in the solutions folder ...
- Beck, A. and L. Tetruashvili. "On the convergence of block coordinate type methods". In: SIAM J. Imaging Sci. 23.4 (2013), pp. 651–694.
- Bottou, L. and O. Bousquet. "The tradeoffs of large scale learning". In: Advances in neural information processing systems. 2008, pp. 161–168.
- Boyd, S. and L. Vandenberghe. Convex optimization. Cambridge University Press, 2004, pp. xiv+716.
- Nocedal, J. "Updating quasi-Newton matrices with limited storage". In: *Mathematics of computation* 35.151 (1980), pp. 773–782.
- Nocedal, J. and S. J. Wright. Numerical optimization. Second. Springer Series in Operations Research and Financial Engineering. New York: Springer, 2006.
- Sun, R. and Y. Ye. "Worst-case complexity of cyclic coordinate descent: O(n²) gap with randomized version". In: Mathematical Programming (2019), pp. 1–34.